

Information Theory and Gradient Descent of Sigmoid Neuron

M medium.com/@manveetdn/notes-on-information-theory-and-gradient-descent-part-of-sigmoid-neuron-padhai-onefourthlabs-8e2f423a62b4

Disclaimer: This is notes on “Information Theory and Gradient Descent Part of Sigmoid Neuron” Lesson (PadhAI onefourthlabs course “A First Course on Deep Learning”)



C	$P(C=?)$	y
	(True) y	(Predicted)
Text	1	0.7
NoText	0	0.3

Text or noText

This is about the Text or noText classification about the true and predicted values.

Expectation:

Lets again, come to the example where we have 4 teams A,B,C,D.

x	A	B	C	D
Based on past Months Probabilities $P(x)$	0.4	0.2	0.1	0.3
Investments $G(x)$	10k	2k	-8k	5k

Here, the expected profit is calculated as

$$(0.4)*(10k)+(0.2)*(2k)+(0.1)*(-8k)+(0.3)*(5k)$$

$$E[p] = \sum_{i \in \{A, B, C, D\}} P(X=i) C(X=i)$$

Formula for Expected profit.

Information Content:

Let us assume **X is a random variable of sun raises in the east**. Here, X is a sure event event it has 4 directions also.

and let **Y be a random variable** of a day with **Storm and noStorm**.

$$IC(X=S) \propto \text{surprise} \propto \frac{1}{P(X=S)}$$

Information content of storm i.e., $X = \text{storm}$ → probability of $X = \text{storm}$

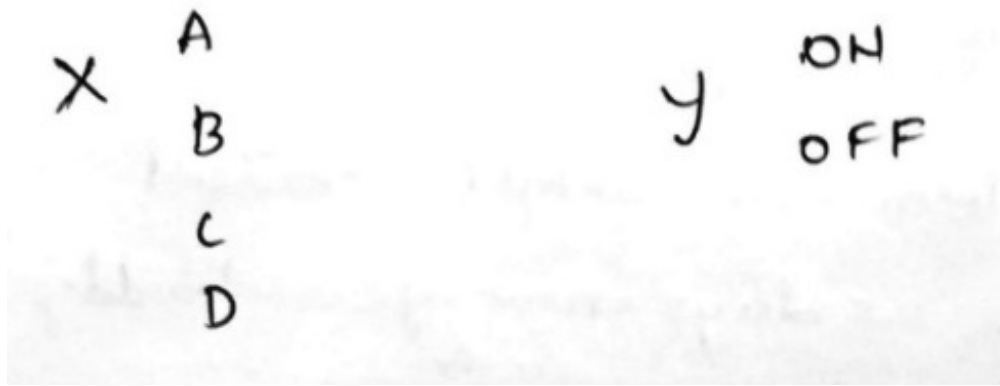
Relation of information content.

Actually **Y = Storm** is a **weighted content because there 365 days in an year hardly very few days will have a storm**. Storm is a rare event and if we tell about that information gained is very high.

Actually, until now Information content is function of the probabilities of the event **IC[P(X=S)]**.

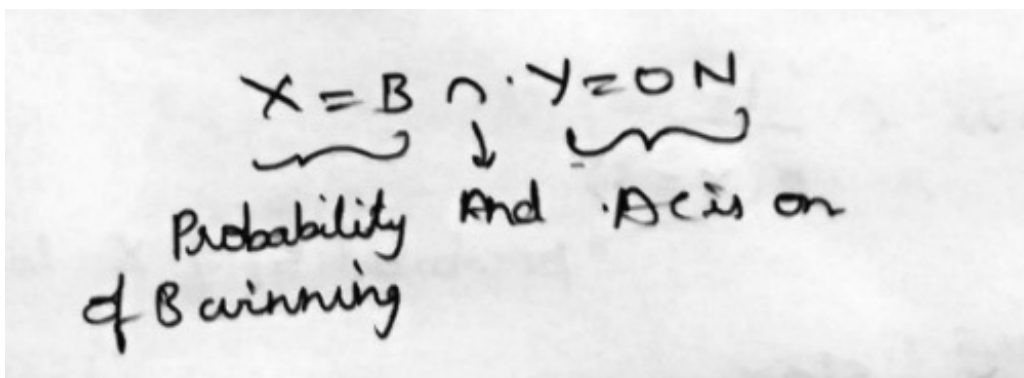
Lets again take some other example some other situation.

There is a Cricket match happening and **X is the Random variable** linked to the probabilities of **winning of team A,B,C,D playing the cricket match**.



Random variables X and Y

Y is a random variable which tells us about the **AC in the room is ON or OFF.**



The above case is about two independent events, like B team winning and AC sos on Condition the Information content of this is

$IC(X=B \cap Y = ON)$ this sum of individual information contents.

$$IC(X=B \cap Y = ON) = IC(X=B) + IC(Y = ON)$$

As we know by our first intuition that Information content is function of the probabilities from above equation we can write it as

$$IC[P(X \cap Y)] = IC[P(X)] + IC[P(Y)]$$

$$IC[P(X).P(Y)] = IC[P(X)] + IC[P(Y)]$$

Which is of the form, **$f(a.b) = f(a) + f(b)$.**

Only **logarithm family of curves** satisfy this form/type of function like **$\log(a.b) = \log(a) + \log(b)$.**

$$IC \propto 1/P(X=A), \text{ fro that } IC = \log(1/P(X=A))$$

$(1/P(X=A))$ comes because more surprise less the information content.

$\log()$ is used to suffice $f(a.b) = f(a) + f(b)$ and only \log satisfies that condition.

$$\begin{aligned}\therefore I C (x=A) &= \log \left(\frac{1}{P(x=A)} \right) \\ &= \log 1 - \log P(x=A) \\ &= -\log P(x=A) \\ \therefore \boxed{I C (x=A) &= -\log_2 P(x=A)}\end{aligned}$$

\log to the base 2 is used as normal practise. Now, we have formula for the Information content.

From this concept we build up entropy and number of bits required to transfer a message.

X	$P(x=?)$	$I C (x=?)$
A	$P(x=A)$	$-\log_2 P(x=A)$
B	$P(x=B)$	$-\log_2 P(x=B)$
C	$P(x=C)$	$-\log_2 P(x=C)$
D	$P(x=D)$	$-\log_2 P(x=D)$

The above table summarises all that we have learned until now and also expectation is given as

$$E[Gain] = \sum_{i \in \{A, B, C, D\}} P(x=i) \text{Gain}(x=i)$$

Expectation is given like this.

Entropy of the random variable:

Entropy of random variable:-

$$H(X) = - \sum_{i \in A, B, C, D} P(X=i) \log_2 P(X=i)$$

It is the expected information content of a Random variable

lets take $P(X=i)$ as P_i

$$\therefore H(X) = - \sum_{i \in \{A, B, C, \dots\}} P_i \log P_i$$

Relation of number of bits:

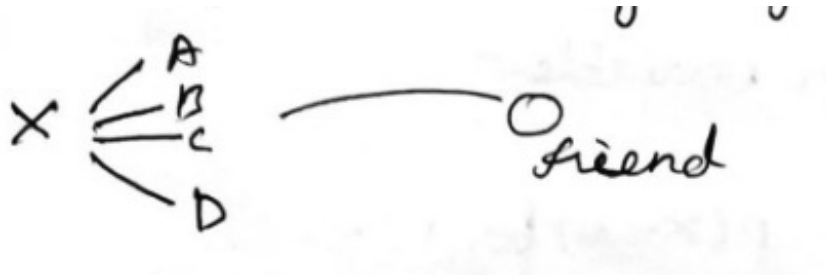
We know

$$H(X) = \sum_{i \in \{A, B, C, \dots\}} P_i \log P_i \quad \text{where } P_i = P(X=i)$$

Lets expand a bit more on this and transmit the no of bits required to transmit a message.

Suppose there is a message X we need to transmit and it can contain 4 values A, B, C, D.

In this 4 command you are going to transmit one.



Transmit a message.

According to Digital logic we need 2 bit to transmit 4 messages 0 0-A, 0 1-B, 1 0-C, 1 1-D.

Here, for every message transmitting you need 2 bits.

			$P(x=?)$	
A	0	0	$\frac{1}{4}$	} Probabilities
B	0	1	$\frac{1}{4}$	
C	1	0	$\frac{1}{4}$	
D	1	1	$\frac{1}{4}$	

Information content of each is given as

$$-\log_2 P(X=?) = -\log_2 \frac{1}{4} = -\log_2 2^{-2} = 2$$

		$P(X=?)$	$-\log_2 P(X=?)$
A	0 0	$\frac{1}{4}$	2
B	0 1	$\frac{1}{4}$	2
C	1 0	$\frac{1}{4}$	2
D	1 1	$\frac{1}{4}$	2

\therefore from this we can say that no. of bits required is equal to the Information content of the message

Let's take other case where message are A B C D
E F G H

Now we will use 3 bits for this

		$P(X=?)$	$-\log_2 P(X=?)$
A	000	$\frac{1}{8}$	3
B	001	$\frac{1}{8}$	3
C	010	$\frac{1}{8}$	3
D	011	$\frac{1}{8}$	3
E	100	$\frac{1}{8}$	3
F	101	$\frac{1}{8}$	3
G	110	$\frac{1}{8}$	3
H	111	$\frac{1}{8}$	3

\therefore Here also the ^{total no. of} number of bit required is equal to the Information content of the message

When you wanna send messages very frequently and the probability of sending each message varies then.

	Equally likely	diff prob	$-\log_2 P$
A	$\frac{1}{4}$	$\frac{1}{2}$	1
B	$\frac{1}{4}$	$\frac{1}{4}$	2
C	$\frac{1}{4}$	$\frac{1}{8}$	3
D	$\frac{1}{4}$	$\frac{1}{8}$	3

frequently sending messages
where probability send A is high
and remaining also differ

A	B	C	D	
$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$P(x=?)$
1	2	3	5	IC

$$\text{Avg}_{\text{bits}} = \frac{1}{2}(1) + \frac{1}{4}(2) + \frac{1}{8}(3) + \frac{1}{8} \times 3$$

$$\text{Avg}_{\text{bits used}} = 1.75$$

\therefore We are going to use less number of bits than
the bit occupied

KL Divergence and Cross Entropy:

$$\begin{array}{c}
 A \quad B \quad C \quad D \\
 x \quad [\quad y_1 \quad y_2 \quad y_3 \quad y_4 \quad] \\
 H(x) \quad [-\log_2 y_1 \quad -\log_2 y_2 \quad -\log_2 y_3 \quad -\log_2 y_4]
 \end{array}$$

$$-\sum y_i \log y_i \quad i = A, B, C, D \text{ and } y \text{ is probability}$$

True Entropy

In the above case the true entropy is given by $(-\sum y_i \log(y_i))$ $i = A, B, C, D$ and y is probability.

Lets consider the n balls earn and when you asked to predict A, B, C, D by peeping into it once and \hat{y} is the one you predict when y is actual value then

$$\begin{array}{c}
 \hat{y} \quad [\quad \hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3 \quad \hat{y}_4 \quad] \cdot \hat{y} \\
 H(\hat{x}) \quad [-\log_2 \hat{y}_1 \quad -\log_2 \hat{y}_2 \quad -\log_2 \hat{y}_3 \quad -\log_2 \hat{y}_4]
 \end{array}$$

Now the actual or expected number of bits you are using will be as below.

Now the actual or expected number of bits you are using will be $-\sum y_i \log \hat{y}_i$ and not $-\sum \hat{y}_i \log \hat{y}_i$

$-\sum y_i \log \hat{y}_i$
 \downarrow
 This the probability
 with which we are going to get them

\downarrow
 value associated with each message
 is $\log_2 y_i$
 because that what we estimated

Now from the above concept we come into

KL Divergence:

If you know the **true distribution** we estimate the average as $-\sum y_i \log(y_i)$

As we don't know the **true distribution** as we estimated it average is $-\sum y_i \log(\hat{y}_i)$

Here average means the average number of bits required

Handwritten notes showing the formulas for cross entropy and entropy. On the left, the formula $-\sum y_i \log \hat{y}_i$ is written, with an arrow pointing from \hat{y}_i to the text $H_{y, \hat{y}}$ (cross entropy). On the right, the formula $-\sum y_i \log y_i$ is written, with an arrow pointing from y_i to the text H_y (Entropy).

Cross entropy and Entropy.

Now we will find the difference between the two distributions.

distance between y and \hat{y} and we can use the difference of bits used in the **cross entropy case and entropy case**.

Handwritten formula for KL divergence as the distance between two distributions. The formula is $D(y || \hat{y}) = -\sum y_i \log \hat{y}_i + \sum y_i \log y_i$. Below the formula, the text "Distance between both" is written.

And that is called KL Divergence

Handwritten formula for KL divergence with labels for cross entropy and entropy. The formula is $KL(y || \hat{y}) = -\sum y_i \log \hat{y}_i + \sum y_i \log y_i$. The first term is labeled "Cross Entropy" and the second term is labeled "Entropy". Below the formula, the text "Computing difference b/w two distributions" is written.

Putting all together:

x
MUMBAI

$$[1 \ 0] y$$

$$[0.2 \ 0.5] \hat{y}$$

$$\hat{y} = \frac{1}{1 + e^{-(wx + b)}}$$

$$\text{Square Error loss, loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$KLD(y || \hat{y}) = - \sum y_i \log \hat{y}_i + \sum y_i \log y_i$$

We need to minimise the above
expression w.r. to w, b

$$\min_{w, b} KLD(y || \hat{y}) = \min_{w, b} \left[- \sum y_i \log \hat{y}_i + \sum y_i \log y_i \right]$$

When we want to make minimisation with respect to the w, b then,

then cross entropy loss $-\sum y_i \log \hat{y}_i$

is a function of w, b and
the normal entropy $\sum y_i \log y_i$ is a constant

$$\min_{(w,b)} f(w,b) + c$$

$$= \min_{(w,b)} f(w,b)$$

\therefore minimization of whole expression minimize
finally as the cross entropy

$$KLD' = \min_{(w,b)} - \sum_{i \in \{Text, NoText\}} y_i \log \hat{y}_i$$

Now we can take a sample data set as an example and solve that using the **KL-Divergence**.

$$\begin{bmatrix} 1 & 0 \end{bmatrix} y$$

$$\begin{bmatrix} 0.7 & 0.3 \end{bmatrix} \hat{y}$$

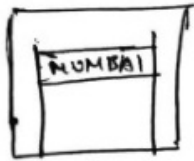
$$\therefore KLD = -1 \log 0.7 - 0 \log 0.3$$

$$KLD = -\log 0.7 = -\log \hat{y}_0$$

$\swarrow \quad \searrow$
 0 1
 NoText Text

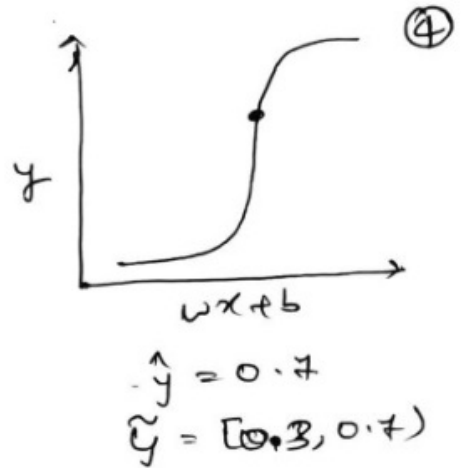
Now ~~now~~ onward we are going used cross
entropy loss function insted of square error loss

Now we are going proceed with two images one having text and another not having text.



$x = \text{image}$
 $y = [0, 1]$

$$\hat{y} = \frac{1}{1 + e^{-(w \cdot x + b)}}$$



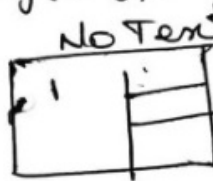
$$L(\theta) = - \sum_{i \in \{0,1\}} y_i \log \tilde{y}_i$$

$$L(\theta) = -(y_0 \log \tilde{y}_0 + y_1 \log \tilde{y}_1)$$

$$L(\theta) = -(y_0 \log(1 - \tilde{y}_1) + y_1 \log \tilde{y}_1)$$



$y = [0, 1]$
 $\hat{y} = 0.7$
 $\tilde{y} = [0.3, 0.7]$



$y = [1, 0]$
 $\hat{y} = 0.2$
 $\tilde{y} = [0.3, 0.2]$

$$L(\theta) = -(y_0 \log(1 - \tilde{y}_1) + y_1 \log \tilde{y}_1) \quad L(\theta) = -(y_0 \log(1 - \tilde{y}_1) + y_1 \log \tilde{y}_1)$$

$$L(\theta) = -(0 * \log 0.3 + 1 * \log 0.7) \quad L(\theta) = -(1 * \log 0.8 + 0 * \log 0.2)$$

$$= -\log 0.7 = -\log \hat{y} \quad = -\log 0.8 = -\log(1 - \hat{y})$$

here $\hat{y} = 0.7$ & $L(\theta) = -\log \hat{y}$

here $\hat{y} = 0.2$ & $L(\theta) = -\log(1 - \hat{y})$

∴ loss function is given as

$$L(\theta) = -((1-y)\log(1-\hat{y}) + y\log\hat{y})$$

When true output is 1 ($y=1$)

$$L(\theta) = -\log\hat{y} \text{ which is same as before}$$

when true output is 0 ($y=0$)

$$L(\theta) = -\log(1-\hat{y}) \text{ which is same as before}$$

$$\therefore L(\theta) = -((1-y)\log(1-\hat{y}) + y\log\hat{y})$$

Learning Algorithm:-

Initialize w, b

Iterate over data

compute \hat{y}

compute $L(w, b)$

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

till satisfied

Input output

1.2	0
-2.1	1
3.2	0
-0.5	1
0.6	1

$$\hat{y} = \frac{1}{1 + e^{-(wx+b)}}$$

$$L(\theta) = -\sum_{i=1}^n ((1-y_i)\log(1-\hat{y}_i) + y_i\log\hat{y}_i)$$

$$L(\theta) = -[(1-y)\log(1-\hat{y}) + y\log\hat{y}]$$

$$\text{We need to compute } \Delta w = \frac{\partial L(\theta)}{\partial w} = \frac{\partial L(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w}$$

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \hat{y}} &= \frac{\partial}{\partial \hat{y}} \{-(1-y)\log(1-\hat{y}) - y\log\hat{y}\} & \frac{\partial \hat{y}}{\partial w} &= \frac{\partial}{\partial w} \left(\frac{1}{1 + e^{-(wx+b)}} \right) \\ &= (-1)(-1) \frac{(1-y)}{(1-\hat{y})} - \frac{y}{\hat{y}} & &= -\frac{1}{(1 + e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)}) \\ &= \frac{\hat{y}(1-y) - y(1-\hat{y})}{(1-\hat{y})\hat{y}} & &= \frac{1}{(1 + e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)}) \\ &= \frac{\hat{y} - y}{(1-\hat{y})\hat{y}} & &= \frac{1}{(1 + e^{-(wx+b)})} \left(\frac{e^{-(wx+b)}}{(1 + e^{-(wx+b)})} \right) * (-1) \\ & & &= \hat{y} * (1-\hat{y}) * x \end{aligned}$$

$$\Delta w = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w} = \frac{\hat{y} - y}{(1-\hat{y})\hat{y}} * \hat{y} * (1-\hat{y}) * x$$

$$\Delta w = (\hat{y} - y) * x$$

like that we will compute Δw and place in the algorithm

```

X = [0.5, 2.5]
Y = [0.2, 0.9]

def f(w, b, x):
    #sigmoid with parameters w, b
    return 1.0 / (1.0 + np.exp(-(w*x + b)))

def error(w, b):
    err = 0.0
    for x, y in zip(X, Y):
        fx = f(w, b, x)
        err += 0.5 * (fx - y) ** 2
    return err

def grad_b(w, b, x, y):
    fx = f(w, b, x)
    return (fx - y) * fx * (1 - fx)

def grad_w(w, b, x, y):
    fx = f(w, b, x)
    return (fx - y) * fx * (1 - fx) * x

def do_gradient_descent():
    w, b, eta = -2, -2, 1.0
    max_epochs = 1000
    for i in range(max_epochs):
        dw, db = 0, 0
        for x, y in zip(X, Y):
            dw += grad_w(w, b, x, y)
            db += grad_b(w, b, x, y)
        w = w - eta * dw
        b = b - eta * db

```

The code of sigmoid neuron before using cross entropy.

```

def grad_w(w, b, x, y):
    fx = f(w, b, x)
    return (fx - y) * x

```

```

def grad_b(w, b, x, y):
    fx = f(w, b, x)
    return (fx - y)

```

```
def error(w, b):
    err = 0.0
    for x, y in zip(X, Y):
        fx = f(w, b, x)
        err += - [(1 - y) * math.log(1 - fx, 2) + y * math.log(fx, 2)]
    return err
```

Changes that need to be made in the error, grad_w, grad_b function.

Small changes in the code need to be done for this cross entropy loss error function, grad_b, grad_w changes will take place

This is all about using gradient descent and cross entropy loss for the sigmoid neuron ⁽⁵⁾

```
def error(w, b):
    err = 0
    for x, y in zip(X, Y):
        fx = f(w, b, x)
        err += - [(1 - y) * math.log(1 - fx, 2) + y * math.log(fx, 2)]
    return err
```

This is the cross entropy loss function

```
def grad_b(w, b, x, y):
    fx = f(w, b, x)
    return (fx - y)
```

This formula which we derived for Δb

```
def grad_w(w, b, x, y):
    fx = f(w, b, x)
    return (fx - y) * x
```

*This is the formula for the Δw which also we derived before $\Delta w = (y - y) * x$*

This is a small try, uploading the notes. I believe in **"Sharing knowledge is that best way of developing skills"**. Comments will be appreciated. Even small edits can be suggested.

Each Applause will be a great encouragement.

Do follow my medium for more updates.....